

Express Mail Label No. EV301223183US
Docket No.: 393032042000
(PATENT)

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re Patent Application of:
Takahiro KAWASHIMA

Application No.: Not Yet Assigned

Filed: Concurrently Herewith

Art Unit: Not Yet Assigned

For: INTERCHANGE FORMAT OF VOICE DATA
IN MUSIC FILE

Examiner: Not Yet Assigned

CLAIM FOR PRIORITY AND SUBMISSION OF DOCUMENT

MS Patent Application
Commissioner for Patents
P.O. Box 1450
Alexandria, VA 22313-1450

Dear Sir:

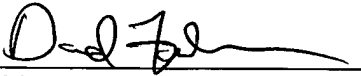
Applicant hereby claims priority under 35 U.S.C. 119 based on the following prior foreign application filed in the following foreign country on the date indicated:

<u>Country</u>	<u>Application No.</u>	<u>Date</u>
Japan	2002-335233	November 19, 2002

In support of this claim, a certified copy of the said original foreign application is filed herewith.

Dated: November 17, 2003

Respectfully submitted,

By 

David L. Fehrman

Registration No.: 28,600
MORRISON & FOERSTER LLP
555 West Fifth Street, Suite 3500
Los Angeles, California 90013
(213) 892-5587

日本国特許庁
JAPAN PATENT OFFICE

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

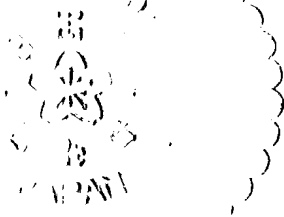
This is to certify that the annexed is a true copy of the following application as filed with this Office.

出願年月日 2002年11月19日
Date of Application:

出願番号 特願2002-335233
Application Number:

[ST. 10/C]: [JP 2002-335233]

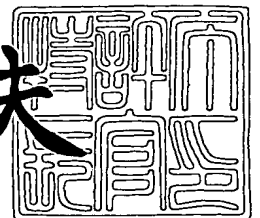
出願人 ヤマハ株式会社
Applicant(s):



2003年 9月22日

特許庁長官
Commissioner,
Japan Patent Office

今井康夫



出証番号 出証特2003-3077641

【書類名】 特許願

【整理番号】 YC30621

【提出日】 平成14年11月19日

【あて先】 特許庁長官殿

【国際特許分類】 G10L 3/00

【発明者】

 【住所又は居所】 静岡県浜松市中沢町 1 0 番 1 号 ヤマハ株式会社内

 【氏名】 川嶋 隆宏

【特許出願人】

 【識別番号】 000004075

 【氏名又は名称】 ヤマハ株式会社

【代理人】

 【識別番号】 100102635

 【弁理士】

 【氏名又は名称】 浅見 保男

【選任した代理人】

 【識別番号】 100106459

 【弁理士】

 【氏名又は名称】 高橋 英生

【選任した代理人】

 【識別番号】 100105500

 【弁理士】

 【氏名又は名称】 武山 吉孝

【選任した代理人】

 【識別番号】 100103735

 【弁理士】

 【氏名又は名称】 鈴木 隆盛

【手数料の表示】

【予納台帳番号】 037338

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9808721

【包括委任状番号】 0106838

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 シーケンスデータのデータ交換フォーマット、音声再生装置及びサーバー装置

【特許請求の範囲】

【請求項 1】 音源を用いて音声を再生させる音声再生シーケンスデータのデータ交換フォーマットであって、

管理用の情報を含むコンテンツデータチャンクと、音声再生シーケンスデータを含むトラックチャンクを有し、

前記音声再生シーケンスデータは、音声の再生を指示する音声再生イベントと、その音声再生イベントを実行するタイミングを先行する音声再生イベントからの経過時間により指定するデュレーションとの組が時間順に配置されているものであることを特徴とするシーケンスデータのデータ交換フォーマット。

【請求項 2】 楽曲を再生させるための楽曲シーケンスデータと音声を再生させるための音声再生シーケンスデータとを含むシーケンスデータのデータ交換フォーマットであって、

前記楽曲シーケンスデータは、演奏イベントデータとその演奏イベントを実行するタイミングを先行する演奏イベントからの経過時間により指定するデュレーションデータとの組が時間順に配置されたものであり、

前記音声再生シーケンスデータは、音声の再生を指示する音声再生イベントデータとその音声再生イベントを実行するタイミングを先行する音声再生イベントからの経過時間により指定するデュレーションデータとの組が時間順に配置されたものであり、

前記楽曲シーケンスデータと前記音声再生シーケンスデータの再生を同時に開始させることにより、当該楽曲と当該音声とを同じ時間軸上で再生することができるようになされていることを特徴とするシーケンスデータのデータ交換フォーマット。

【請求項 3】 前記楽曲シーケンスデータと前記音声再生シーケンスデータはそれぞれ異なるチャンクに含まれていることを特徴とする請求項 2 記載のシーケンスデータのデータ交換フォーマット。

【請求項4】 前記音声再生イベントデータは、（１）合成される音声の読みを示すテキスト情報と音声表現を指定する韻律記号とからなるテキスト記述型の情報、（２）合成される音声を示す音素情報と韻律制御情報とからなる音素記述型の情報、又は、（３）再生される音声を示すフレーム時間毎のフォルマント制御情報からなるフォルマントフレーム記述型の情報、の再生を指示するデータであることを特徴とする請求項1乃至3のいずれかに記載のシーケンスデータのデータ交換フォーマット。

【請求項5】 音声の再生を指示する音声再生イベントデータと、その音声再生イベントを実行するタイミングを先行する音声再生イベントからの経過時間により指定するデュレーションデータとの組により構成されているシーケンスデータに基づいて音声を再生する音声再生装置であって、

テキスト情報及び韻律記号と音素情報及び韻律制御情報との対応を記録した辞書を参照して、テキスト情報及び韻律記号に対応する音素情報及び韻律制御情報を取得する第1の手段と、

音素情報及び韻律制御情報とフォルマント制御情報との対応を記録した辞書を参照して、音素情報及び韻律制御情報に対応するフォルマント制御情報を取得する第2の手段と、

フォルマント制御情報をフレーム毎に音源部に出力する出力手段とを有し、

前記音声再生イベントデータが合成される音声の読みを示すテキスト情報と韻律記号とからなるテキスト記述型の情報の再生を指示するデータであるときは、前記第1の手段、前記第2の手段及び前記出力手段を用いて当該音声を再生し、

前記音声再生イベントデータが合成される音声に対応する音素情報と韻律制御情報を含む音素記述型の情報の再生を指示するデータであるときは、前記第2の手段と前記出力手段を用いて当該音声を再生し、

前記音声再生イベントデータがフレーム毎のフォルマント制御情報の再生を指示するデータであるときは、前記出力手段を用いて当該音声を再生する

ことを特徴とする音声再生装置。

【請求項6】 音声の再生を指示する音声再生イベントデータと、その音声再生イベントを実行するタイミングを先行する音声再生イベントからの経過時間

により指定するデュレーションデータとの組により構成されているシーケンスデータに基づいて音声を再生する音声再生装置であって、

音素情報及び韻律制御情報とフォルマント制御情報との対応を記憶した辞書を参照して、音素情報及び韻律制御情報に対応するフォルマント制御情報を取得する手段と、

フォルマント制御情報をフレーム毎に音源部に出力する出力手段とを有し、

前記音声再生イベントデータが合成される音声に対応する音素情報と韻律制御情報を含む音素記述型の情報の再生を指示するデータであるときは、前記取得する手段と前記出力手段とを用いて当該音声を再生し、

前記音声再生イベントデータがフレーム毎のフォルマント制御情報の再生を指示するデータであるときは、前記出力手段を用いて当該音声を再生する

ことを特徴とする音声再生装置。

【請求項 7】 音声の再生を指示する音声再生イベントデータと、その音声再生イベントを実行するタイミングを先行する音声再生イベントからの経過時間により指定するデュレーションデータとの組により構成されているシーケンスデータに基づいて音声を再生する音声再生装置であって、

フォルマント制御情報をフレーム毎に音源部に出力する出力手段を有し、

前記音声再生イベントデータがフレーム毎のフォルマント制御情報の再生を指示するデータであるときに、前記出力手段を用いて当該音声を再生することを特徴とする音声再生装置。

【請求項 8】 シーケンスデータを蓄積し、クライアント装置からの要求に応じて対応するシーケンスデータを配信するサーバー装置であって、

前記シーケンスデータは、

楽曲シーケンスデータを含むチャンクと音声再生シーケンスデータを含むチャンクを有し、

前記楽曲シーケンスデータは、演奏イベントデータとその演奏イベントを実行するタイミングを先行する演奏イベントからの経過時間により指定するデュレーションデータとの組が時間順に配置されたものであり、

前記音声再生シーケンスデータは、音声の再生を指示する音声再生イベントデ

ータとその音声再生イベントを実行するタイミングを先行する音声再生イベントからの経過時間により指定するデュレーションデータとの組が時間順に配置されたものである

ことを特徴とするサーバー装置。

【請求項 9】 前記音声再生イベントデータは、（１）合成される音声の読みを示すテキスト情報と韻律記号とからなるテキスト記述型の情報、（２）合成される音声を示す音素情報と韻律制御情報とからなる音素記述型の情報、又は、（３）再生される音声を示すフレーム時間毎のフォルマント制御情報からなるフォルマントフレーム記述型の情報、の再生を指示するデータであり、

クライアント装置からの要求に応じて、前記（１）～（３）のうちのいずれかの音声再生イベントデータを含む音声再生シーケンスデータを選択して、当該クライアント装置に配信することを特徴とする請求項 8 記載のサーバー装置。

【発明の詳細な説明】

【 0 0 0 1 】

【発明の属する技術分野】

本発明は、シーケンスデータのデータ交換フォーマット、音声再生装置及びサーバー装置に関する。

【 0 0 0 2 】

【従来の技術】

音源を用いて音楽を表現するためのデータを頒布したり、相互に利用したりするためのデータ交換フォーマットとして、SMF（Standard MIDI file format）やSMAF（Synthetic Music Mobile Application Format）などが知られている。SMAFは、携帯端末などにおいてマルチメディアコンテンツを表現するためのデータフォーマット仕様である（非特許文献 1 参照）。

【 0 0 0 3 】

図 1 5 を参照しつつ SMAF について説明する。

この図において、1 0 0 は SMAF ファイルであり、チャンクとよばれるデータの塊が基本構造となっている。チャンクは固定長（8 バイト）のヘッダ部と任意長のボディ部とからなり、ヘッダ部は、さらに、4 バイトのチャンク ID と 4

バイトのチャンクサイズに分けられる。チャンク ID はチャンクの識別子に用い、チャンクサイズはボディ部の長さを示している。SMAF ファイルは、それ自体及びそれに含まれる各種データも全てチャンク構造となっている。

この図に示すように、SMAF ファイル 1 0 0 の中身は、管理用の情報が格納されているコンテンツ・インフォ・チャンク (Contents Info Chunk) 1 0 1 と、出力デバイスに対するシーケンスデータを含む 1 つ以上のトラックチャンク 1 0 2 ~ 1 0 8 とからなる。シーケンスデータは出力デバイスに対する制御を時間を追って定義したデータ表現である。1 つの SMAF ファイル 1 0 0 に含まれる全てのシーケンスデータは時刻 0 で同時に再生を開始するものと定義されており、結果的に全てのシーケンスデータが同期して再生される。

シーケンスデータはイベントとデュレーションの組み合わせで表現される。イベントは、シーケンスデータに対応する出力デバイスに対する制御内容のデータ表現であり、デュレーションは、イベントとイベントとの間の経過時間を表現するデータである。イベントの処理時間は実際には 0 ではないが、SMAF のデータ表現としては 0 とみなし、時間の流れは全てデュレーションで表わすようにしている。あるイベントを実行する時刻は、そのシーケンスデータの先頭からのデュレーションを積算することで一意に決定することができる。イベントの処理時間は、次のイベントの処理開始時刻に影響しないことが原則である。従って、値が 0 のデュレーションを挟んで連続したイベントは同時に実行すると解釈される。

【0 0 0 4】

SMAF では、前記出力デバイスとして、MIDI (musical instrument digital interface) 相当の制御データで発音を行う音源デバイス 1 1 1、PCM データの再生を行う PCM 音源デバイス (PCM デコーダ) 1 1 2、テキストや画像の表示を行う LCD などの表示デバイス 1 1 3 などが定義されている。

トラックチャンクには、定義されている各出力デバイスに対応して、スコアトラックチャンク 1 0 2 ~ 1 0 5、PCM オーディオトラックチャンク 1 0 6、グラフィックストラックチャンク 1 0 7 及びマスタートラックチャンク 1 0 8 がある。ここで、マスタートラックチャンクを除くスコアトラックチャンク、PCM

オーディオトラックチャンク及びグラフィックストラックチャンクは、それぞれ最大 2 5 6 トラックまで記述することが可能である。

図示する例では、スコアトラックチャンク 1 0 2 ～ 1 0 5 は音源デバイス 1 1 1 を再生するためのシーケンスデータを格納し、PCMトラックチャンク 1 0 6 は PCM 音源デバイス 1 1 2 で発音される ADPCM や MP3、TwinVQ 等の wave データをイベント形式で格納し、グラフィックトラックチャンク 1 0 7 は背景画や差込静止画、テキストデータと、それらを表示デバイス 1 1 3 で再生するためのシーケンスデータを格納している。また、マスタートラックチャンク 1 0 8 には SMAF シーケンサ自身を制御するためのシーケンスデータが格納されている。

【0 0 0 5】

一方、音声合成の手法として、LPC などのフィルタ合成方式や複合正弦波音声合成法などの波形合成方式がよく知られている。複合正弦波音声合成法（CSM 法）は、複数の正弦波の和により音声信号をモデル化し音声合成を行う方式であり、簡単な合成法でありながら良質な音声を合成することができる。（非特許文献 2 参照）。

また、音源を用いて音声合成させることにより、歌声を発生させる音声合成装置も提案されている（特許文献 1 参照）。

【0 0 0 6】

【非特許文献 1】

SMAF 仕様書 Ver. 3.06 ヤマハ株式会社、[平成 1 4 年 1 0 月 1 8 日検索]、インターネット<URL: <http://smaf.yamaha.co.jp>>

【非特許文献 2】

嵯峨山茂樹、板倉文忠、「複合正弦波音声合成方式の検討と合成器の試作」、日本音響学会、音声研究会資料、資料番号 S80-12(1980-5)、p.93-100、(1980.5.26)

【特許文献 1】

特開平 9 - 5 0 2 8 7 号公報

【0 0 0 7】

【発明が解決しようとする課題】

上述のように、SMAFは、MIDI相当のデータ（楽曲データ）、PCMオーディオデータ、テキストや画像の表示データなどの各種シーケンスデータを含み、全シーケンスを時間的に同期して再生することができる。

しかしながら、SMFやSMAFには音声（人の声）を表現することについては、定義されていない。

そこで、SMFなどのMIDIイベントを拡張して音声を合成することも考えられるが、この場合は、音声部分のみ一括して取り出して音声合成するときに処理が複雑になるという問題点がある。

【0008】

そこで本発明は、柔軟性があり、かつ、楽曲シーケンスなどと音声再生シーケンスとを同期して再生させることが可能なシーケンスデータのデータ交換フォーマット、該データ交換フォーマットのファイルを再生することができる音声再生装置及び該データ交換フォーマットのデータを配信することができるサーバー装置を提供することを目的としている。

【0009】

【課題を解決するための手段】

上記目的を達成するために、本発明のシーケンスデータのデータ交換フォーマットは、音源を用いて音声を再生させる音声再生シーケンスデータのデータ交換フォーマットであって、管理用の情報を含むコンテンツデータチャンクと、音声再生シーケンスデータを含むトラックチャンクを有し、前記音声再生シーケンスデータは、音声の再生を指示する音声再生イベントと、その音声再生イベントを実行するタイミングを先行する音声再生イベントからの経過時間により指定するデュレーションとの組が時間順に配置されているものである。

また、本発明の他のシーケンスデータのデータ交換フォーマットは、楽曲を再生させるための楽曲シーケンスデータと音声を再生させるための音声再生シーケンスデータとを含むシーケンスデータのデータ交換フォーマットであって、前記楽曲シーケンスデータは、演奏イベントデータとその演奏イベントを実行するタイミングを先行する演奏イベントからの経過時間により指定するデュレーションデータとの組が時間順に配置されたものであり、前記音声再生シーケンスデータ

は、音声の再生を指示する音声再生イベントデータとその音声再生イベントを実行するタイミングを先行する音声再生イベントからの経過時間により指定するデュレーションデータとの組が時間順に配置されたものであり、前記楽曲シーケンスデータと前記音声再生シーケンスデータの再生を同時に開始させることにより、当該楽曲と当該音声とを同じ時間軸上で再生することができるようになされているものである。

さらに、前記楽曲シーケンスデータと前記音声再生シーケンスデータはそれぞれ異なるチャンクに含まれているものである。

そして、前記音声再生イベントデータは、(1) 合成される音声の読みを示すテキスト情報と音声表現を指定する韻律記号とからなるテキスト記述型の情報、(2) 合成される音声を示す音素情報と韻律制御情報とからなる音素記述型の情報、又は、(3) 再生される音声を示すフレーム時間毎のフォルマント制御情報からなるフォルマントフレーム記述型の情報、の再生を指示するデータとされているものである。

【0010】

さらにまた、本発明の音声再生装置は、音声の再生を指示する音声再生イベントデータと、その音声再生イベントを実行するタイミングを先行する音声再生イベントからの経過時間により指定するデュレーションデータとの組により構成されているシーケンスデータに基づいて音声を再生する音声再生装置であって、テキスト情報及び韻律記号と音素情報及び韻律制御情報との対応を記録した辞書を参照して、テキスト情報及び韻律記号に対応する音素情報及び韻律制御情報を取得する第1の手段と、音素情報及び韻律制御情報とフォルマント制御情報との対応を記録した辞書を参照して、音素情報及び韻律制御情報に対応するフォルマント制御情報を取得する第2の手段と、フォルマント制御情報をフレーム毎に音源部に出力する出力手段とを有し、前記音声再生イベントデータが合成される音声の読みを示すテキスト情報と韻律記号とからなるテキスト記述型の情報の再生を指示するデータであるときは、前記第1の手段、前記第2の手段及び前記出力手段を用いて当該音声を再生し、前記音声再生イベントデータが合成される音声に対応する音素情報と韻律制御情報を含む音素記述型の情報の再生を指示するデー

タであるときは、前記第2の手段と前記出力手段を用いて当該音声を再生し、前記音声再生イベントデータがフレーム毎のフォルマント制御情報の再生を指示するデータであるときは、前記出力手段を用いて当該音声を再生するものである。

又は、前記第2の手段と前記出力手段を有するもの、もしくは、前記出力手段を有するものである。

【0011】

さらにまた、本発明のサーバー装置は、シーケンスデータを蓄積し、クライアント装置からの要求に応じて対応するシーケンスデータを配信するサーバー装置であって、前記シーケンスデータは、楽曲シーケンスデータを含むチャンクと音声再生シーケンスデータを含むチャンクを有し、前記楽曲シーケンスデータは、演奏イベントデータとその演奏イベントを実行するタイミングを先行する演奏イベントからの経過時間により指定するデュレーションデータとの組が時間順に配置されたものであり、前記音声再生シーケンスデータは、音声の再生を指示する音声再生イベントデータとその音声再生イベントを実行するタイミングを先行する音声再生イベントからの経過時間により指定するデュレーションデータとの組が時間順に配置されたものとされている。

そして、前記音声再生イベントデータは、(1) 合成される音声の読みを示すテキスト情報と韻律記号とからなるテキスト記述型の情報、(2) 合成される音声を示す音素情報と韻律制御情報とからなる音素記述型の情報、又は、(3) 再生される音声を示すフレーム時間毎のフォルマント制御情報からなるフォルマントフレーム記述型の情報、の再生を指示するデータであり、クライアント装置からの要求に応じて、前記(1)～(3)のうちのいずれかの音声再生イベントデータを含む音声再生シーケンスデータを選択して、当該クライアント装置に配信するようになされている。

【0012】

【発明の実施の形態】

図1は、本発明における音声再生シーケンスデータのデータ交換フォーマットの一実施の形態を示す図である。この図において、1は本発明のデータ交換フォーマットを有するファイルである。このファイル1は、前述したSMAFファイ

ルと同様に、チャンク構造を基本としており、ヘッダ部とボディ部とを有する（ファイルチャンク）。

前記ヘッダ部には、ファイルを識別するためのファイルID（チャンクID）と後続するボディ部の長さを示すチャンクサイズが含まれている。

ボディ部はチャンク列であり、図示する例では、コンテンツ・インフォ・チャンク（Contents Info Chunk）2、オプション・データ・チャンク（Optional Data Chunk）3、及び、音声再生シーケンスデータを含むHV（Human Voice）トラックチャンク4が含まれている。なお、図1には、HVトラックチャンク4として、HVトラックチャンク#00の一つのみが記載されているが、ファイル1中に複数個のHVトラックチャンク4を含ませることができる。

また、本発明においては、前記HVトラックチャンク4に含まれる音声再生シーケンスデータとして、3つのフォーマットタイプ（TSeq型、PSeq型、FSeq型）が定義されている。これらについては後述する。

前記コンテンツ・インフォ・チャンク2には、含まれているコンテンツのクラス、種類、著作権情報、ジャンル名、曲名、アーティスト名、作詞/作曲者名などの管理用の情報が格納されている。また、前記著作権情報やジャンル名、曲名、アーティスト名、作詞/作曲者名などの情報を格納するオプション・データ・チャンク3を設けても良い。

【0013】

図1に示した音声再生シーケンスデータのデータ交換フォーマットは、それ単独で音声を再生することができるが、前記HVトラックチャンク4をデータチャンクの一つとして前述したSMAFファイルに含ませることができる。

図2は、上述したHVトラックチャンク4をデータチャンクの一つとして含む本発明のシーケンスデータのデータ交換フォーマットを有するファイルの構造を示す図である。このファイルは、SMAFファイルを音声再生シーケンスデータを含むように拡張したものであるということが出来る。なお、この図において、前記図15に示したSMAFファイル100と同一の要素には同一の番号を付す。

この図に示すように、前述した音声再生シーケンスデータのデータ交換フォー

マットにおけるHVトラックチャंक4を、前述したスコアトラックチャंक102～105、PCMオーディオトラックチャंक106、グラフィックストラックチャंक107などと共に、SMAFファイル100中に格納することにより、楽曲の演奏や画像、テキストの表示と同期して音声を再生することが可能となり、例えば、楽音に対し、音源が歌うコンテンツなどを実現することができるようになる。

【0014】

図3は、前記図2に示した本発明のデータ交換フォーマットのファイルを作成するシステム及び該データ交換フォーマットファイルを利用するシステムの概略構成の一例を示す図である。

この図において、21はSMFやSMAFなどの楽曲データファイル、22は再生される音声に対応するテキストファイル、23は本発明によるデータ交換フォーマットのファイルを作成するためのデータ・フォーマット制作ツール（オーサリング・ツール）、24は本発明のデータ交換フォーマットを有するファイルである。

オーサリング・ツール23は、再生する音声の読みを示す音声合成用テキストファイル22を入力して、編集作業などを行い、それに対応する音声再生シーケンスデータを作成する。そして、SMFやSMAFなどの楽曲データファイル21に該作成した音声再生シーケンスデータを加えて、本発明のデータ交換フォーマット仕様に基づくファイル（前記図2に示したHVトラックチャंकを含むSMAFファイル）24を作成する。

【0015】

作成されたファイル24は、シーケンスデータに含まれているデュレーションにより規定されるタイミングで音源部27に制御パラメータを供給するシーケンサ26と、シーケンサ26から供給される制御パラメータに基づいて音声を再生出力する音源部27を有する利用装置25に転送され、そこで、楽曲などとともに音声が同期して再生されることとなる。

図4は前記音源部27の概略構成の一例を示す図である。

この図に示した例では、音源部27は、複数のフォルマント生成部28と1個

のピッチ生成部 29 を有しており、前記シーケンサ 26 から出力されるフォルマント制御情報（各フォルマントを生成するためのフォルマント周波数、レベルなどのパラメータ）及びピッチ情報に基づいて各フォルマント生成部 28 で対応するフォルマント信号を発生し、これらをミキシング部 30 で加算することにより対応する音声合成出力が生成される。なお、各フォルマント生成部 28 はフォルマント信号を発生させるためにその元となる基本波形を発生させるが、この基本波形の発生には、例えば、周知の FM 音源の波形発生器を利用することができる。

【0016】

前述のように、本発明においては、前記 H V トラックチャック 4 に含まれる音声再生シーケンスデータに 3 つのフォーマットタイプを用意し、これらを任意に選択して用いることができるようにしている。以下、これらについて説明する。

再生する音声を記述するためには、再生する音声に対応する文字情報、言語に依存しない発音情報、音声波形そのものを示す情報など抽象度が異なる各種の段階の記述方法があるが、本発明においては、（a）テキスト記述型（TSeq型）、（b）音素記述型（PSeq型）及び（c）フォルマント・フレーム記述型（FSeq型）の 3 通りのフォーマットタイプを定義している。

【0017】

まず、図 5 を参照して、これら 3 つのフォーマットタイプの相違について説明する。

（a）テキスト記述型（TSeq型）

TSeq型は、発音すべき音声をテキスト表記により記述するフォーマットであり、それぞれの言語による文字コード（テキスト情報）とアクセントなどの音声表現を指示する記号（韻律記号）とを含む。このフォーマットのデータはエディタなどを用いて直接作成することができる。再生するときは、図 5 の（a）に示すように、ミドルウェア処理により、該 TSeq型のシーケンスデータを、まず、PSeq型に変換し（第 1 のコンバート処理）、次に、PSeq型をFSeq型に変換（第 2 のコンバート処理）して、前記音源部 27 に出力することとなる。

ここで、TSeq型からPSeq型へ変換する第 1 のコンバート処理は、言語に依存す

る情報である文字コード（例えば、ひらがなやカタカナなどのテキスト情報）と韻律記号と、それに対応する言語に依存しない発音を示す情報（音素）と韻律を制御するための韻律制御情報を格納した第1の辞書を参照することにより行われ、PSeq型からFSeq型への変換である第2のコンバート処理は、各音素及び韻律制御情報とそれに対応するフォルマント制御情報（各フォルマントを生成するためのフォルマントの周波数、帯域幅、レベルなどのパラメータ）を格納した第2の辞書を参照することにより行われる。

（b）音素記述型（PSeq型）

PSeq型は、SMFで定義するMIDIイベントに類似する形式で発音すべき音声に関する情報を記述するものであり、音声記述としては言語依存によらない音素単位をベースとする。図5の（b）に示すように、前記オーサリング・ツールなどを用いて実行されるデータ制作処理においては、まずTSeq型のデータファイルを作成し、これを第1のコンバート処理によりPSeq型に変換する。このPSeq型を再生するときは、ミドルウェア処理として実行される第2のコンバート処理によりPSeq型のデータファイルをFSeq型に変換して、音源部27に出力する。

（c）フォルマント・フレーム記述型（FSeq型）

FSeq型は、フォルマント制御情報をフレーム・データ列として表現したフォーマットである。図5の（c）に示すように、データ制作処理において、TSeq型→第1のコンバート処理→PSeq型→第2のコンバート処理→FSeq型への変換を行う。また、サンプリングされた波形データから通常の音声分析処理と同様の処理である第3のコンバート処理によりFSeq型のデータを作成することもできる。再生時には、該FSeq型のファイルをそのまま前記音源部に出力して再生することができる。

このように、本発明においては、抽象度の異なる3種類のフォーマットタイプを定義し、個々の場合に依じて、所望のタイプを選択することができるようにしている。また、音声を再生するために実行する前記第1のコンバート処理及び前記第2のコンバート処理をミドルウェア処理として実行させることにより、アプリケーションの負担を軽減することができる。

【0018】

次に、前記H Vトラックチャンク 4（図 1）の内容について詳細に説明する。

前記図 1 に示したように、各H Vトラックチャンク 4 には、このH Vトラックチャンクに含まれている音声再生シーケンスデータが前述した 3 通りのフォーマットタイプのうちのどのタイプであるかを示すフォーマットタイプ（Format Type）、使用されている言語種別を示す言語タイプ（Language Type）及びタイムベース（Timebase）をそれぞれ指定するデータが記述されている。

フォーマットタイプ（Format Type）の例を表 1 に示す。

【表 1】

フォーマットタイプ	説 明
0x00	TSeq型
0x01	PSeq型
0x02	FSeq型

【 0 0 1 9】

言語タイプ（Language Type）の例を表 2 に示す。

【表 2】

言語タイプ	説 明
0x00	Shift-JIS
0x02	EUC-KR(KS)

なお、ここでは、日本語（0x00；0xは 1 6 進を表わす。以下、同じ。）と韓国語（0x01）のみを示しているが、中国語、英語などその他の言語についても同様に定義することができる。

【 0 0 2 0】

タイムベース（Timebase）は、このトラックチャンクに含まれるシーケンスデータチャンク内のデュレーション及びゲートタイムの基準時間を定めるものである。この実施の形態では、20msecとされているが任意の値に設定することができる。

る。

【表 3】

タイムベース	説 明
0x11	20 msec

【0 0 2 1】

前述した 3 通りのフォーマットタイプのデータの詳細についてさらに説明する

。

(a) Tseq型 (フォーマットタイプ=0x00)

前述のように、このフォーマットタイプは、テキスト表記によるシーケンス表現 (TSeq: text sequence) を用いたフォーマットであり、シーケンスデータチャンク 5 と n 個 (n は 1 以上の整数) の TSeq データチャンク (TSeq#00～TSeq#n) 6, 7, 8 を含んでいる (図 1)。シーケンスデータに含まれる音声再生イベント (ノートオンイベント) で TSeq データチャンクに含まれるデータの再生を指示する。

【0 0 2 2】

(a-1) シーケンスデータチャンク

シーケンスデータチャンクは、SMAF におけるシーケンスデータチャンクと同様に、デュレーションとイベントの組み合わせを時間順に配置したシーケンスデータを含む。図 6 の (a) はシーケンスデータの構成を示す図である。ここで、デュレーションは、イベントとイベントの間の時間を示している。先頭のデュレーション (Duration 1) は、時刻 0 からの経過時間を示している。図 6 の (b) は、イベントがノートメッセージである場合に、デュレーションとノートメッセージに含まれるゲートタイムの関係を示す図である。この図に示すように、ゲートタイムはそのノートメッセージの発音時間を示している。なお、図 6 で示したシーケンスデータチャンクの構造は、PSeq 型及びFSeq 型におけるシーケンスデータチャンクにおいても同様である。

このシーケンスデータチャンクでサポートされるイベントとしては、次の 3 通

りのイベントがある。なお、以下に記述する初期値は、イベント指定がないときのデフォルト値である。

(a-1-1) ノートメッセージ「0x9n kk gt」

ここで、n：チャンネル番号（0x0[固定]）、kk：TSeqデータ番号（0x00～0x7F）、gt：ゲートタイム（1～3 バイト）である。

ノートメッセージは、チャンネル番号 n で指定されるチャンネルの TSeq データ番号 kk で指定される TSeq データチャンクを解釈し発音を開始するメッセージである。なお、ゲートタイム gt が「0」のノート・メッセージについては発音を行わない。

(a-1-2) ボリューム「0xBn 0x07 vv」

ここで、n：チャンネル番号（0x0[固定]）、vv：コントロール値（0x00～0x7F）である。なお、チャンネルボリュームの初期値は0x64である。

ボリュームは、指定チャンネルの音量を指定するメッセージである。

(a-1-3) パン「0xBn 0x0A vv」

ここで、n：チャンネル番号（0x0[固定]）、v v：コントロール値（0x00～0x7F）である。なお、パンポット初期値は、0x40（センター）である。

パンメッセージは、指定チャンネルのステレオ音場位置を指定するメッセージである。

【 0 0 2 3 】

(a-2) TSeqデータチャンク（TSeq#00～TSeq#n）

TSeqデータチャンクは、音声合成用の情報として、言語や文字コードに関する情報、発音する音の設定、（合成する）読み情報を表記したテキストなどを含んだ、しゃべり用フォーマットでありタグ形式で書かれている。このTSeqデータチャンクは、ユーザーによる入力を容易にするためテキスト入力となっている。

タグは、“<”（0x3C）で始まり制御タグと値が続く形式であり、TSeqデータチャンクはタグの列で構成されている。ただし、スペースは含まず、制御タグ及び値に“<”は使用することはできない。また、制御タグは必ず1文字とする。制御タグとその有効値に例を下の表4に示す。

【 0 0 2 4 】

【表 4】

タグ		値	意味
L	(0x4C)	Language	言語情報
C	(0x43)	<i>code</i>	文字コード名
T	(0x54)	全角文字列	合成用テキスト
P	(0x50)	0-	無音の挿入
S	(0x53)	0-127	再生速度
V	(0x56)	0-127	音量
N	(0x4E)	0-127	音の高さ
G	(0x47)	0-127	音色選択
R	(0x52)	None	リセット
Q	(0x51)	None	終了

【0 0 2 5】

前記制御タグのうちのテキストタグ「T」について、さらに説明する。

テキストタグ「T」に後続する値は、全角ひらがな文字列で記述された読み情報（日本語の場合）と音声表現を指示する韻律記号（Shift-JISコード）からなる。文末にセンテンス区切り記号がないときは、“。”で終わるのと同じ意味とする。

以下に示すのは韻律記号であり、読み情報の文字の後につく。

”、”(0x8141)：センテンスの区切り（通常のイントネーション）。

”。”(0x8142)：センテンスの区切り（通常のイントネーション）。

”？”(0x8148)：センテンスの区切り（疑問のイントネーション）。

”””(0x8166)：ピッチを上げるアクセント（変化後の値はセンテンス区切りまで有効）。

”__”(0x8151)：ピッチを下げるアクセント（変化後の値はセンテンス区切りまで有効）。

”ー”(0x815B)：長音（直前の語を長く発音する。複数でより長くなる。）

【0 0 2 6】

図 7 の（a）は、TSeqデータチャンクのデータの一例を示す図であり、（b）

はその再生時間処理について説明するための図である。

最初のタグ「<LJAPANESE」で言語が日本語であることを示し、「<CS-JIS」で文字コードがシフト J I S であること、「<G4」で音色選択（プログラムチェンジ）、「<V1000」で音量の設定、「<N64」で音の高さを指定している。「<T」は合成用テキストを示し、「<P」はその値により規定される msec 単位の無音期間の挿入を示している。

図 7 の（b）に示すように、この TSeq データチャンクのデータは、デュレーションにより指定されるスタート時点から 1000 msec の無音期間をおいた後に、「い' やーー、き__よーわ' さ__むい__ねー。」と発音され、その後 1500 msec の無音期間をおいた後に「こ' のままい__ったら、は' ちが__つわ、た' いへ' ん__やねー。」と発音される。ここで、「'」、「__」、「-」に応じてそれぞれに対応するアクセントや長音の制御が行われる。

【 0 0 2 7 】

このように、TSeq 型は、各国語それぞれに特化した発音をするための文字コードと音声表現（アクセントなど）をタグ形式で記述したフォーマットであるため、エディタなどを用いて直接作成することができる。従って、TSeq データチャンクのファイルはテキストベースで容易に加工することができ、例えば、記述されている文章からイントネーションを変更したり、語尾を加工することで方言に対応するといったことを容易に行うことができる。また、文章中の特定単語だけを入れ替えることも容易にできる。さらに、データ・サイズが小さいという長所がある。

一方、この TSeq 型データチャンクのデータを解釈し音声合成をするための処理負荷が大きくなる、より細かいピッチ制御ができにくい、フォーマットを拡張し複雑な定義を増やせば、ユーザ・フレンドリーでなくなってしまう、言語（文字）コードに依存する（例えば、日本語の場合には Shift-JIS が一般であるが、他国語の場合には、それに応じた文字コードでフォーマットを定義する必要がある。）などという短所がある。

【 0 0 2 8 】

（b）PSeq 型（フォーマットタイプ=0x01）

このPSeq型は、M I D I イベントに類似する形式の音素によるシーケンス表現（PSeq：phoneme sequence）を用いたフォーマットタイプである。この形式は、音素を記述するようにしているので言語依存がない。音素は発音を示す文字情報により表現することができ、例えば、複数の言語に共通にアスキーコードを用いることができる。

前記図 1 に示したように、このPSeq型は、セットアップ・データ・チャンク 9、ディクショナリ・データ・チャンク 1 0 及びシーケンス・データ・チャンク 1 1 を含んでいる。シーケンスデータ中の音声再生イベント（ノートメッセージ）で指定されたチャンネルの音素と韻律制御情報の再生を指示する。

【0 0 2 9】

(b-1) セットアップ・データ・チャンク（Setup Data Chunk）（オプション）

音源部分の音色データなどを格納するチャンクであり、イクスクルーシブ・メッセージの並びを格納する。この実施の形態では、含まれているイクスクルーシブ・メッセージは、H V 音色パラメータ登録メッセージである。

H V 音色パラメータ登録メッセージは「0xF0 Size 0x43 0x79 0x07 0x7F 0x01 PC data ... 0xF7」というフォーマットであり、PC：プログラム番号（0x02～0x0F）、data：H V 音色パラメータである。

このメッセージは、該当するプログラム番号PCのH V 音色パラメータを登録する。

【0 0 3 0】

H V 音色パラメータを次の表 5 に示す。



【表 5】

#0	基本音声番号
#1	ピッチシフト量[Cent]
#2	フォルマント周波数シフト量 1
#3	フォルマント周波数シフト量 2
#4	:
#5	フォルマント周波数シフト量 n
#6	フォルマントレベルシフト量 1
#7	フォルマントレベルシフト量 2
#8	:
#9	フォルマントレベルシフト量 n
#10	オペレータ波形選択 1
#11	オペレータ波形選択 2
#12	:
#13	オペレータ波形選択 n

【0 0 3 1】

表 5 に示すように、H V 音色パラメータとしては、ピッチシフト量、第 1 ～ 第 n（n は 2 以上の整数）の各フォルマントに対するフォルマント周波数シフト量、フォルマントレベルシフト量及びオペレータ波形選択情報が含まれている。前述のように、処理装置内には、各音素とそれに対応するフォルマント制御情報（フォルマントの周波数、帯域幅、レベルなど）を記述したプリセット辞書（第 2 の辞書）が記憶されており、H V 音色パラメータは、このプリセット辞書に記憶されているパラメータに対するシフト量を規定している。これにより、全ての音素について同様のシフトが行われ、合成される音声の声質を変化させることができる。

なお、この H V 音色パラメータにより、0x02～0x0F に対応する数（すなわち、プログラム番号の数）の音色を登録することができる。

【0 0 3 2】

(b-2) ディクショナリデータチャンク (Dictionary Data Chunk) (オプション)

このチャンクには、言語種別に応じた辞書データ、例えば、前記プリセット辞書と比較した差分データやプリセット辞書で定義していない音素データなどを含む。

む辞書データを格納する。これにより、音色の異なる個性のある音声を合成することが可能となる。

【 0 0 3 3 】

(b-3) シーケンスデータチャンク (Sequence Data Chunk)

前述のシーケンスデータチャンクと同様に、デュレーションとイベントの組み合わせを時間順に配置したシーケンスデータを含む。

このPSeq型におけるシーケンスデータチャンクでサポートするイベント（メッセージ）を次に列挙する。読み込み側は、これらのメッセージ以外は無視する。また、以下に記述する初期設定値は、イベント指定がないときのデフォルト値である。

【 0 0 3 4 】

(b-3-1) ノートメッセージ「0x9n Nt Vel Gatetime Size data ...」

ここで、n：チャンネル番号（0x0[固定]）、Nt：ノート番号（絶対値ノート指定：0x00～0x7F，相対値ノート指定：0x80～0xFF）、Vel：ベロシティ（0x00～0x7F）、Gatetime：ゲートタイム長（Variable）、Size：データ部のサイズ（可変長）である。

このノートメッセージにより、指定チャンネルの音声の発音が始まる。

なお、ノート番号のMSBは、解釈を絶対値と相対値とに切り替えるフラグである。MSB以外の7ビットはノート番号を示す。音声の発音はモノラルのみであるため、ゲートタイムが重なる場合は後着優先として発音する。オーサリング・ツールなどでは、重なりのあるデータは作られないように制限を設けることが望ましい。

【 0 0 3 5 】

データ部は、音素とそれに対する韻律制御情報（ピッチベンド、ボリューム）を含み、次の表6に示すデータ構造からなる。

【表 6】

#0	ディレイ
#1	音素数 [=n]
#2	音素 1
#3	:
#4	音素 n
#5	音素ピッチベンド数 [=N]
#6	音素ピッチベンド位置 1
#7	音素ピッチベンド 1
#8	:
#9	音素ピッチベンド位置 N
#10	音素ピッチベンド N
#11	音素ボリューム数 [=M]
#12	音素ボリューム位置 1
#13	音素ボリューム 1
#14	:
#15	音素ボリューム位置 M
#16	音素ボリューム M

【0036】

表 6 に示すように、データ部は、音素の数 n（#1）、例えばアスキーコードで記述した個々の音素（音素 1～音素 n）（#2～#4）、及び、韻律制御情報からなっている。韻律制御情報はピッチベンドとボリュームであり、ピッチベンドに関して、その発音区間を音素ピッチベンド数（#5）により規定される N 個の区間に区切り、それぞれにおけるピッチベンドを指定するピッチベンド情報（音素ピッチベンド位置 1，音素ピッチベンド 1（#6～#7）～音素ピッチベンド位置 N，音素ピッチベンド N（#9～#10））と、ボリュームに関して、その発音区間を音素ボリューム数（#11）により規定される M 個の区間に区切り、それぞれにおけるボリュームを指定するボリューム情報（音素ボリューム位置 1，音素ボリューム 1（#12, #13）～音素ボリューム位置 M，音素ボリューム M（#15, #16））からなっている。

【0037】

図 8 は、前記韻律制御情報について説明するための図である。ここでは、発音する文字情報が「o h a y o u」である場合を例にとって示している。また、こ

の例では、 $N=M=128$ としている。この図に示すように、発音する文字情報（「o h a y o u」）に対応する区間を $128 (=N=M)$ の区間に区切り、各点におけるピッチとボリュームを前記ピッチベンド情報及びボリューム情報で表現して韻律を制御するようにしている。

【 0 0 3 8 】

図 9 は、前記ゲートタイム長 (Gatetime) とディレイタイム (Delay Time (#0)) との関係を示す図である。この図に示すように、ディレイタイムにより、実際の発音をデューレーションで規定されるタイミングよりも遅らせることができる。なお、Gate time = 0 は、禁止とする。

【 0 0 3 9 】

(b-3-2) プログラムチェンジ 「0xCn pp」

ここで、n：チャンネル番号 (0x0[固定])、pp：プログラム番号 (0x00～0xF F) である。また、プログラム番号の初期値は 0x00 とされている。

このプログラムチェンジメッセージにより指定されたチャンネルの音色が設定される。ここで、チャンネル番号は、0x00：男声プリセット音色、0x01：女声プリセット音色、0x02～0x0F：拡張音色である。

【 0 0 4 0 】

(b-3-3) コントロールチェンジ

コントロールチェンジメッセージとしては、次のものがある。

(b-3-3-1) チャンネルボリューム 「0xBn 0x07 vv」

ここで、n：チャンネル番号 (0x0[固定])、vv：コントロール値 (0x00～0x7F) である。また、チャンネルボリュームの初期値は 0x64 とされている。

このチャンネルボリュームメッセージは、指定チャンネルの音量を指定するものであり、チャンネル間の音量バランスを設定することを目的としている。

(b-3-3-2) パン 「0xBn 0x0A vv」

ここで、n：チャンネル番号 (0x0[固定])、vv：コントロール値 (0x00～0x7F) である。パンポットの初期値は 0x40 (センター) とされている。

このメッセージは、指定チャンネルのステレオ音場位置を指定する。

【 0 0 4 1 】

(b-3-3-3) エクスプレッション 「0xBn 0x0B vv」

ここで、n：チャンネル番号 (0x0[固定])、vv：コントロール値 (0x00～0x7F) である。このエクスプレッションメッセージの初期値は0x7F (最大値) とされている。

このメッセージは、指定チャンネルのチャンネル・ボリュームで設定した音量の変化を指定する。これは曲中で音量を変化させる目的で使用される。

【 0 0 4 2 】

(b-3-3-4) ピッチベンド 「0xEn ll mm」

ここで、n：チャンネル番号 (0x0[固定])、ll：ベンド値 L S B (0x00～0x7F)、mm：ベンド値 M S B (0x00～0x7F) である。ピッチベンドの初期値は M S B 0x40、L S B 0x00 とされている。

このメッセージは、指定チャンネルのピッチを上下に変化させる。変化幅 (ピッチ・ベンド・レンジ) の初期値は ± 2 半音であり、0x00 / 0x00 で下方向へのピッチ・ベンドが最大となる。0x7F / 0x7F で上方向へのピッチ・ベンドが最大となる。

【 0 0 4 3 】

(b-3-3-5) ピッチベンド・センシティビティ 「0x8n bb」

ここで、n：チャンネル番号 (0x0[固定])、bb：データ値 (0x00～0x18) である。このピッチベンド・センシティビティの初期値は0x02である。

このメッセージは、指定チャンネルのピッチ・ベンドの感度設定を行う。単位は半音である。例えば、bb=01 のときは ± 1 半音 (変化範囲は計 2 半音) となる。

【 0 0 4 4 】

このように、PSeq型のフォーマットタイプは、発音を示す文字情報で表現した音素単位をベースとし、M I D I イベントに類似する形式で音声情報を記述したものであり、データ・サイズはTSeq型よりは大きい、FSeq型よりは小さくなる。

これにより、M I D I と同様に時間軸上の細かいピッチやボリュームをコントロールすることができる、音素ベースで記述しているため言語依存性がない、音色 (声質) を細かく編集することができる、M I D I と類似した制御ができ、従

来のM I D I 機器へ追加実装し易いという長所を有している。

一方、文章や単語レベルの加工ができない、処理側において、TSeq型よりは軽いものの、フォーマットを解釈し音声合成するための処理負荷がかかるという短所を有している。

【0 0 4 5】

(c) フォルマント・フレーム記述 (FSeq) 型 (フォーマットタイプ=0x02)

フォルマント制御情報 (各フォルマントを生成するための、フォルマント周波数やゲインなどのパラメータ) をフレーム・データ列として表現したフォーマットである。すなわち、一定時間 (フレーム) の間は、発音する音声のフォルマントなどは一定であるとし、各フレーム毎に発音する音声に対応するフォルマント制御情報 (各々のフォルマント周波数やゲインなど) を更新するシーケンス表現 (FSeq: formant sequence) を用いる。シーケンスデータに含まれるノートメッセージにより指定されたFSeqデータチャンクのデータの再生を指示する。

このフォーマットタイプは、シーケンスデータチャンクと n 個 (n は以上の整数) のFSeqデータチャンク (FSeq#00~FSeq#n) を含んでいる。

【0 0 4 6】

(c-1) シーケンスデータチャンク

前述のシーケンスデータチャンクと同様に、デュレーションとイベントの組を時間順に配置したシーケンスデータを含む。

以下に、このシーケンスデータチャンクでサポートするイベント (メッセージ) を列挙する。読み込み側は、これらのメッセージ以外は無視する。また、以下に記述する初期設定値は、イベント指定がないときのデフォルト値である。

(c-1-1) ノート・メッセージ 「0x9n kk gt」

ここで、n: チャンネル番号 (0x0 [固定])、kk: FSeqデータ番号 (0x00~0x7F)、gt: ゲートタイム (1~3バイト) である。

このメッセージは、指定チャンネルのFSeqデータ番号のFSeqデータチャンクを解釈し発音を開始するメッセージである。なお、ゲートタイムが"0"のノート・メッセージは発音を行わない。

【0 0 4 7】

(c-1-2) ボリューム 「0xBn 0x07 vv」

ここで、n：チャンネル番号 (0x0[固定])、vv：コントロール値 (0x00～0x7F) である。なお、チャンネルボリュームの初期値は0x64である。

このメッセージは、指定チャンネルの音量を指定するメッセージである。

【 0 0 4 8 】

(c-1-3) パン 「0xBn 0x0A vv」

ここで、n：チャンネル番号 (0x0[固定])、vv：コントロール値 (0x00～0x7F) である。なお、パンポットの初期値は0x40 (センター) である。

このメッセージは、指定チャンネルのステレオ音場位置を指定するメッセージである。

【 0 0 4 9 】

(c-2) FSeqデータチャンク (FSeq#00～FSeq#n)

FSeqデータチャンクは、FSeqフレーム・データ列で構成する。すなわち、音声情報を所定時間長 (例えば、20msec) を有するフレーム毎に切り出し、それぞれのフレーム期間内の音声データを分析して得られたフォルマント制御情報 (フォルマント周波数やゲインなど) を、それぞれのフレームの音声データを表わすフレーム・データ列として表現したフォーマットである。

表 7 にFSeqのフレーム・データ列を示す。

【 0 0 5 0 】

【表 7】

#0	オペレータ波形 1
#1	オペレータ波形 2
#2	:
#3	オペレータ波形 n
#4	フォルマントレベル 1
#5	フォルマントレベル 2
#6	:
#7	フォルマントレベル n
#8	フォルマント周波数 1
#9	フォルマント周波数 2
#10	:
#11	フォルマント周波数 n
#12	有声／無声 切り替え

【0051】

表 7 において、#0～#3は音声合成に用いる複数個（この実施の形態においては、n 個）のフォルマントの波形の種類（サイン波、矩形波など）を指定するデータである。#4～#11は、フォルマントレベル（振幅）（#4～#7）と中心周波数（#8～#11）により n 個のフォルマントを規定するパラメータである。#4と#8が第 1 フォルマント（#0）を規定するパラメータ、以下同様に、#5～#7と#9～#11は第 2 フォルマント（#1）～第 n フォルマント（#3）を規定するパラメータである。また、#12は無声／有声を示すフラグなどである。

図 10 は、フォルマントのレベルと中心周波数を示す図であり、この実施の形態においては、第 1～第 n フォルマントまでの n 個のフォルマントのデータを用いるようにしている。前記図 4 に示したように、各フレーム毎の第 1～第 n フォルマントに関するパラメータとピッチ周波数に関するパラメータは、前記音源部 27 のフォルマント生成部とピッチ生成部に供給され、そのフレームの音声合成出力が前述のようにして生成出力される。

【0052】

図 11 は、前記FSeqデータチャンクのボディ部のデータを示す図である。前記表 7 に示したFSeqのフレームデータ列のうち、#0～#3は、各フォルマントの波形

の種類を指定するデータであり、各フレームごとに指定する必要はない。従って、図 1 1 に示すように、最初のフレームについては、前記表 7 に示した全てのデータとし、後続するフレームについては、前記表 7 における #4 以降のデータだけでよい。FSeq データチャンクのボディ部を図 1 1 のようにすることにより、総データ数を少なくすることができる。

【 0 0 5 3 】

このように、FSeq 型は、フォルマント制御情報（各々のフォルマント周波数やゲインなど）をフレーム・データ列として表現したフォーマットであるため、FSeq 型のファイルをそのまま音源部に出力することにより音声を再生することができる。従って、処理側は音声合成処理の必要がなく、CPU は所定時間ごとにフレームを更新する処理を行うのみでよい。なお、既に格納されている発音データに対し、一定のオフセットを与えることで音色（声質）を変更することができる。

ただし、FSeq 型のデータは文章や単語レベルの加工がしづらく、音色（声質）を細かく編集したり、時間軸上の発音長やフォルマント変位を変更することができない。さらに、時間軸上のピッチやボリュームを制御することはできるが、元のデータのオフセットで制御することとなるため、制御しにくいのに加え、処理負荷が増大するという短所がある。

【 0 0 5 4 】

次に、上述したシーケンスデータのデータ交換フォーマットを有するファイルを利用するシステムについて説明する。

図 1 2 は、上述した音声再生シーケンスデータを再生する音声再生装置の一つである携帯通信端末に対し、上述したデータ交換フォーマットのファイルを配信するコンテンツデータ配信システムの概略構成を示す図である。

この図において、5 1 は携帯通信端末、5 2 は基地局、5 3 は前記複数の基地局を統括する移動交換局、5 4 は複数の移動交換局を管理するとともに公衆網などの固定網やインターネット 5 5 とのゲートウェイとなる関門局、5 6 はインターネット 5 5 に接続されたダウンロードセンターのサーバーコンピュータである。

コンテンツデータ制作会社 57 は、前記図 3 に関して説明したように、専用のオーサリング・ツールなどを用い、SMF や SMAF などの楽曲データ及び音声合成用テキストファイルから本発明のデータ交換フォーマットを有するファイルを作成し、サーバーコンピュータ 56 に転送する。

サーバーコンピュータ 56 には、コンテンツデータ制作会社 57 により制作された本発明のデータ交換フォーマットを有するファイル（前記 H V トラックチャックを含む SMAF ファイルなど）が蓄積されており、携帯通信端末 51 や図示しないコンピュータなどからアクセスするユーザーからのリクエストに応じて、対応する前記音声再生シーケンスデータを含む楽曲データなどを配信する。

【0055】

図 13 は、音声再生装置の一例である前記携帯通信端末 51 の一構成例を示すブロック図である。

この図において、61 はこの装置全体の制御を行う中央処理装置（CPU）、62 は各種通信制御プログラムや楽曲再生のためのプログラムなどの制御プログラムおよび各種定数データなどが格納されている ROM、63 はワークエリアとして使用されるとともに楽曲ファイルや各種アプリケーションプログラムなどを記憶する RAM、64 は液晶表示装置（LCD）などからなる表示部、65 はバイブレータ、66 は複数の操作ボタンなどを有する入力部、67 は変復調部などからなりアンテナ 68 に接続される通信部である。

また、69 は、送話マイク及び受話スピーカに接続され、通話のための音声信号の符号化および復号を行う機能を有する音声処理部、70 は前記 RAM 63 などに記憶された楽曲ファイルに基づいて楽曲を再生するとともに、音声を再生して、スピーカ 71 に出力する音源部、72 は前記各構成要素間のデータ転送を行うためのバスである。

ユーザーは、前記携帯通信端末 51 を用いて、前記図 12 に示したダウンロードセンターのサーバー 56 にアクセスし、前記 3 つのフォーマットタイプのうちの所望のタイプの音声再生シーケンスデータを含む本発明のデータ交換フォーマットのファイルをダウンロードして前記 RAM 63 などに格納し、そのまま再生したり、あるいは、着信メロディとして使用することができる。

【 0 0 5 6 】

図 1 4 は、前記サーバーコンピュータ 5 6 からダウンロードして前記 R A M 6 3 に記憶した本発明のデータ交換フォーマットのファイルを再生する処理の流れを示すフローチャートである。ここでは、ダウンロードしたファイルが、前記図 2 に示したフォーマットにおいて、スコアトラックチャンクと H V トラックチャンクを有するファイルであるとして説明する。

楽曲の再生の開始指示があったとき、或いは、着信メロディとして使用する場合は着信が発生して処理が開始されると、ダウンロードしたファイルに含まれている音声部（H V トラックチャンク）と楽曲部（スコアトラックチャンク）を分離する（ステップ S 1）。そして、音声部については、そのフォーマットタイプが（a）TSeq型であるときには、TSeq型をPSeq型に変換する第 1 のコンバート処理とPSeq型をFSeq型に変換する第 2 のコンバート処理を実行してFSeq型に変換し、（b）PSeq型であるときには、前記第 2 のコンバート処理を行ってFSeq型に変換し、（c）FSeq型であるときにはそのままというように、フォーマットタイプに応じた処理を行ってFSeq型のデータに変換し（ステップ S 2）、各フレームのフォルマント制御データをフレーム毎に更新して前記音源部 7 0 に供給する（ステップ S 3）。一方、楽曲部については、音源部に所定のタイミングで楽音発生パラメータを供給する（ステップ S 4）。これにより、音声と楽曲が合成して（ステップ S 5）、出力される（ステップ S 6）。

【 0 0 5 7 】

前記図 3 に関して説明したように、本発明のデータ交換フォーマットは、S M F や S M A F などの既存の楽曲データ 2 1 に音声合成用テキストデータ 2 2 に基づいて作成した音声再生シーケンスデータを付け加えることにより制作することができるため、上述のように着信メロディなどに利用した場合に多種のエンターテイメント性のあるサービスを提供することが可能となる。

【 0 0 5 8 】

また、上記においてはダウンロードセンターのサーバーコンピュータ 5 6 からダウンロードした音声再生シーケンスデータを再生するものであったが、音声再生装置で上述した本発明のデータ交換フォーマットのファイルを作成することも

できる。

前記携帯通信端末 5 1 において、発声したいテキストに対応する前記TSeq型のTSeqデータチャンクを入力部 6 6 から入力する。例えば、「<Tお' っはよー、げ__んき？」と入力する。そして、これをそのまま、あるいは、前記第 1、第 2 のコンバート処理を行って、前述の 3 つのフォーマットタイプのうちのいずれかの音声再生シーケンスデータとし、本発明のデータ交換フォーマットのファイルへ変換して保存する。そして、そのファイルをメールに添付して相手端末に送信する。

このメールを受信した相手方の携帯通信端末では、受信したファイルのタイプを解釈し、対応した処理を行ってその音源部を用いて当該音声再生する。

このように、携帯通信端末で、データを送信する前に加工することで、多種のエンターテインメント性のあるサービスを提供することが可能となる。この場合、それぞれの加工方法で、サービスに最適な音声合成用フォーマット種類を選択する。

【 0 0 5 9 】

さらにまた、近年では、携帯通信端末において J a v a (TM) によるアプリケーションプログラムをダウンロードして実行することができるようになっている。そこで、J a v a (TM) アプリケーションプログラムを用いてより多彩な処理を行わせることができる。

すなわち、携帯通信端末上で、発声したいテキストを入力する。そして、J a v a (TM) アプリケーションプログラムにより、入力されたテキストデータを受け取り、該テキストに合致した画像データ（例えば、しゃべっている顔）を貼付け、本発明のデータ交換フォーマットのファイル（H Vトラックチャンクとグラフィックストラックチャンクを有するファイル）へ変換し、J a v a (TM) アプリケーションプログラムから A P I 経由で本ファイルをミドルウェア（シーケンサ、音源や画像を制御するソフトウェアモジュール）に送信する。ミドルウェアは送られたファイル・フォーマットを解釈し、音源で音声を再生しながら表示部で画像を同期して表示する。

このように、J a v a (TM) アプリケーションのプログラミングにより、多種の

エンターテインメント性のあるサービスを提供することができる。この場合、それぞれの加工方法で、サービスに最適な音声合成用フォーマット種類を選択する。

【0 0 6 0】

なお、上述した実施の形態においては、H Vトラックチャンクに含まれる音声再生シーケンスデータのフォーマットを3つの型に応じて異なるフォーマットとしていたが、これに限られることはない。例えば、前記図1に示したように、(a) TSeq型と(c) FSeq型は、いずれも、シーケンスデータチャンクとTSeqあるいはFSeqデータチャンクを有するものであり、基本的な構造は同一であるので、これらを統一し、データチャンクのレベルで、TSeq型のデータチャンクであるのかFSeq型のデータチャンクであるのかを識別するようにしてもよい。

また、上述した各表に記載したデータの定義は、何れも一例に過ぎないものであり、任意に変更することができる。

【0 0 6 1】

【発明の効果】

以上説明したように、本発明の音声再生シーケンスデータのデータ交換フォーマットによれば、音声再生のためのシーケンスを表現することができるとともに、異なるシステムや装置の間で音声再生シーケンスデータを頒布したり交換することが可能となる。

また、楽曲シーケンスデータと音声再生シーケンスデータを各々異なるチャンクに含むようにした本発明のシーケンスデータのデータ交換フォーマットによれば、1つのフォーマット・ファイルで音声再生シーケンスと楽曲シーケンスの同期を取って再生することができる。

また、楽曲シーケンスデータと音声再生シーケンスデータを独立に記述することができ、一方のみを取り出して再生させることが容易にできる。

また、3つのフォーマットタイプを選択することができる本発明のデータ交換フォーマットによれば、音声再生の用途や処理側の負荷を考慮し、最も適切なフォーマットタイプを選択することができる。

【図面の簡単な説明】

【図1】 本発明における音声再生シーケンスデータのデータ交換フォーマット

ットの一実施の形態を示す図である。

【図 2】 H Vトラックチャンクをデータチャンクの一つとして含む S M A F ファイルの例を示す図である。

【図 3】 本発明のデータ交換フォーマットを作成するシステム及び該データ交換フォーマットファイルを利用するシステムの概略構成の一例を示す図である。

【図 4】 音源部の概略構成の一例を示す図である。

【図 5】 (a) TSeq型、(b) PSeq型、及び、(c) FSeq型の 3 通りのフォーマットタイプの違いについて説明するための図である。

【図 6】 (a) はシーケンスデータの構成、(b) はデュレーションとゲートタイムの関係を示す図である。

【図 7】 (a) はTSeqデータチャンクの一例を示す図であり、(b) はその再生時間処理について説明するための図である。

【図 8】 韻律制御情報について説明するための図である。

【図 9】 ゲートタイムとディレイタイムとの関係を示す図である。

【図 1 0】 フォルマントのレベルと中心周波数を示す図である。

【図 1 1】 FSeqデータチャンクのボディ部のデータを示す図である。

【図 1 2】 音声再生装置の一つである携帯通信端末に対し本発明のデータ交換フォーマットのファイルを配信するコンテンツデータ配信システムの概略構成の一例を示す図である。

【図 1 3】 携帯通信端末の一構成例を示すブロック図である。

【図 1 4】 本発明のデータ交換フォーマットのファイルを再生する処理の流れを示すフローチャートである。

【図 1 5】 S M A F の概念を説明するための図である。

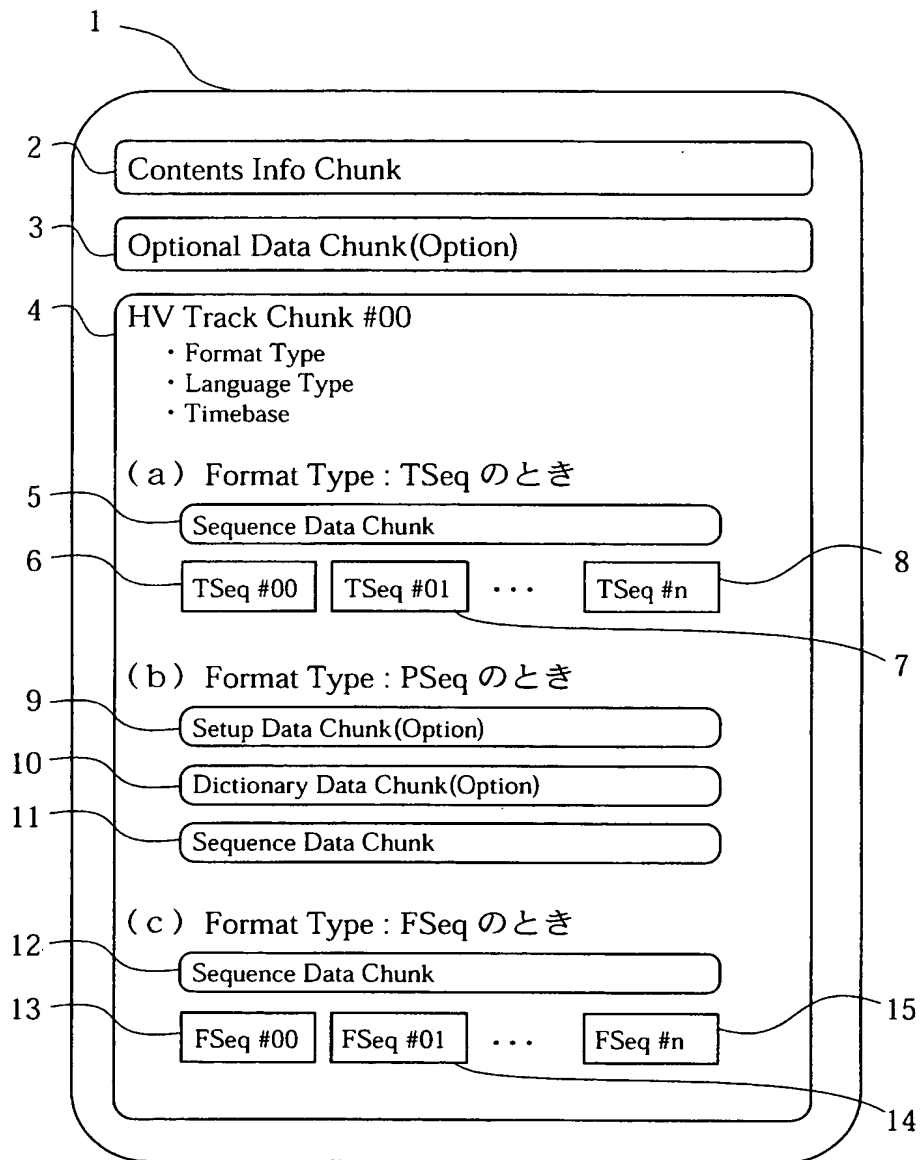
【符号の説明】

1 本発明のデータ交換フォーマットを有するファイル、2 コンテンツ・インフォ・チャンク、3 オプショナル・データ・チャンク、4 H Vトラックチャンク、5, 1 1, 1 2 シーケンスデータチャンク、6 ~ 8 TSeqデータチャンク、9 セットアップデータチャンク、1 0 デクシヨナリデータチャンク

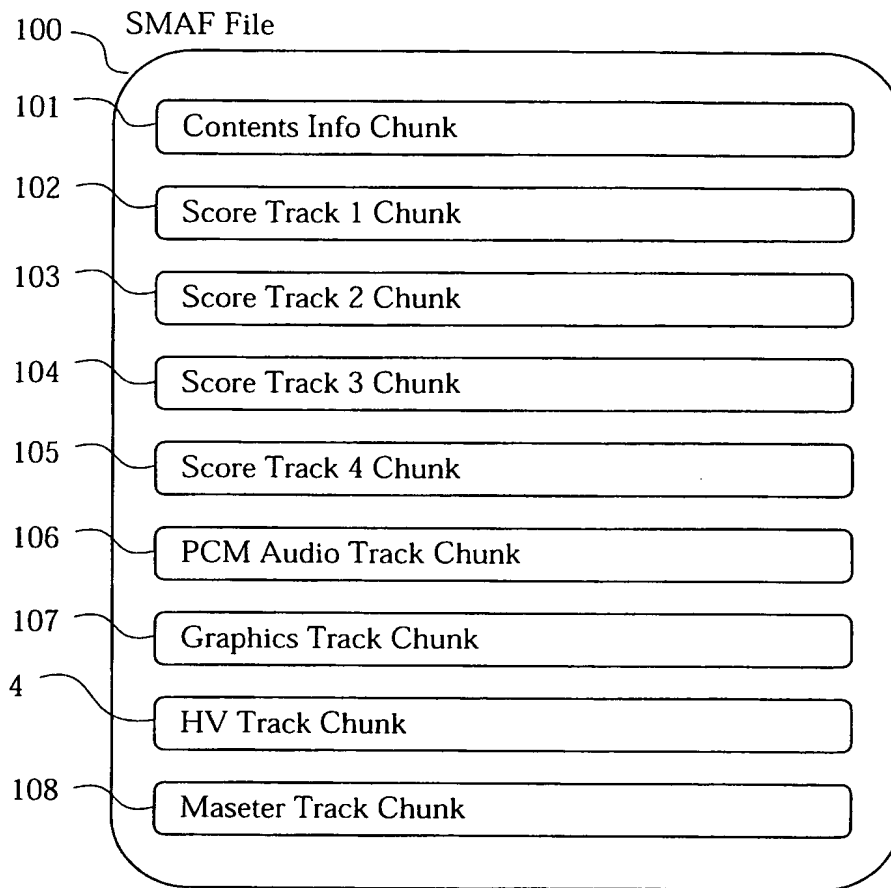
、13～15 FSeqデータチャンク、21 楽曲データ、22 テキストファイル、23 オーサリング・ツール、24 本発明のデータ交換フォーマットを有するファイル、25 利用装置、26 シーケンサ、27 音源部、28 フォルマント生成部、29 ピッチ生成部、30 ミキシング部、51 携帯通信端末、52 基地局、53 移動交換局、54 関門局、55 インターネット、56 ダウンロードサーバー、57 コンテンツデータ制作会社

【書類名】 図面

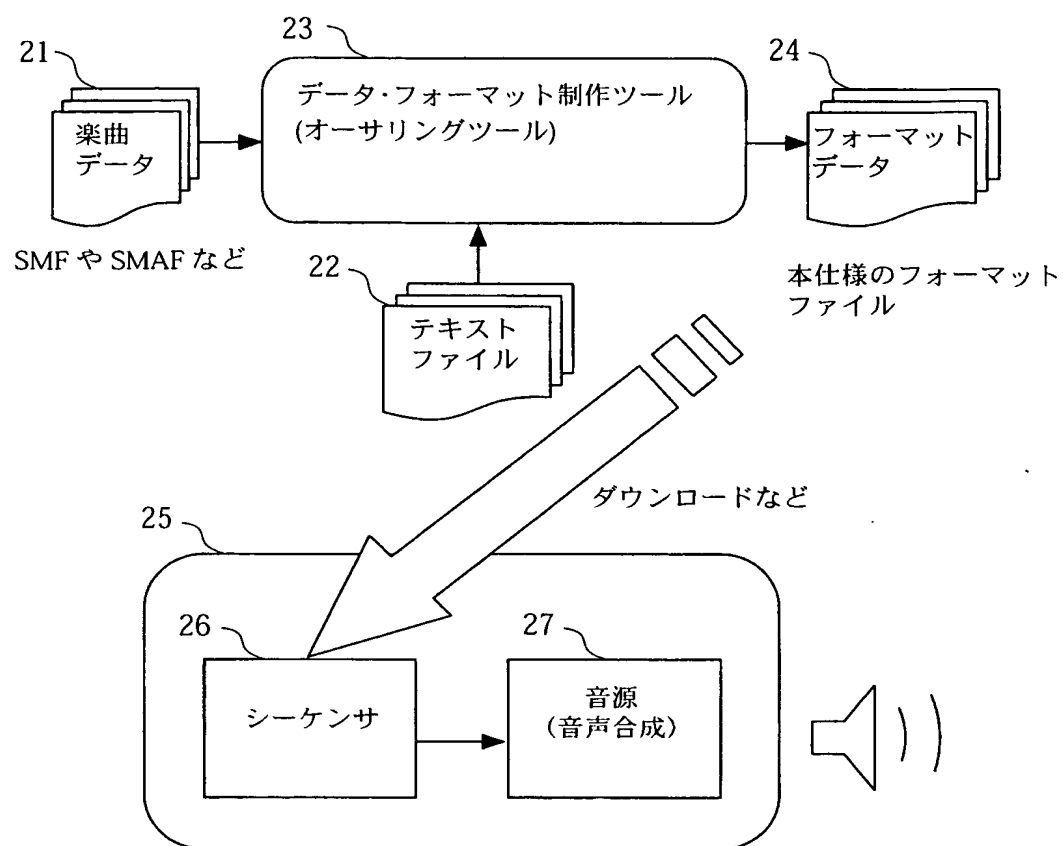
【図 1】



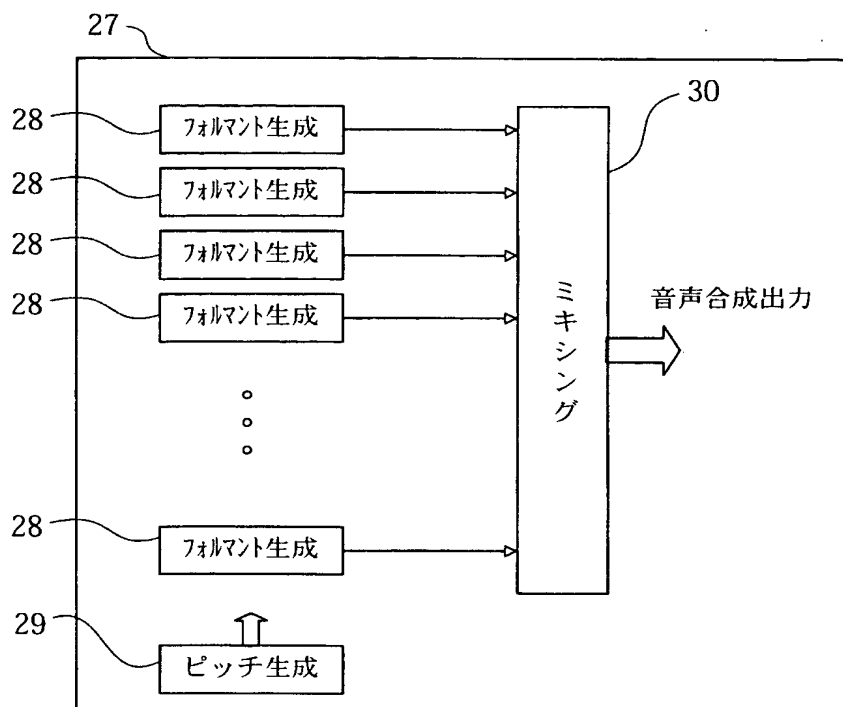
【図 2】



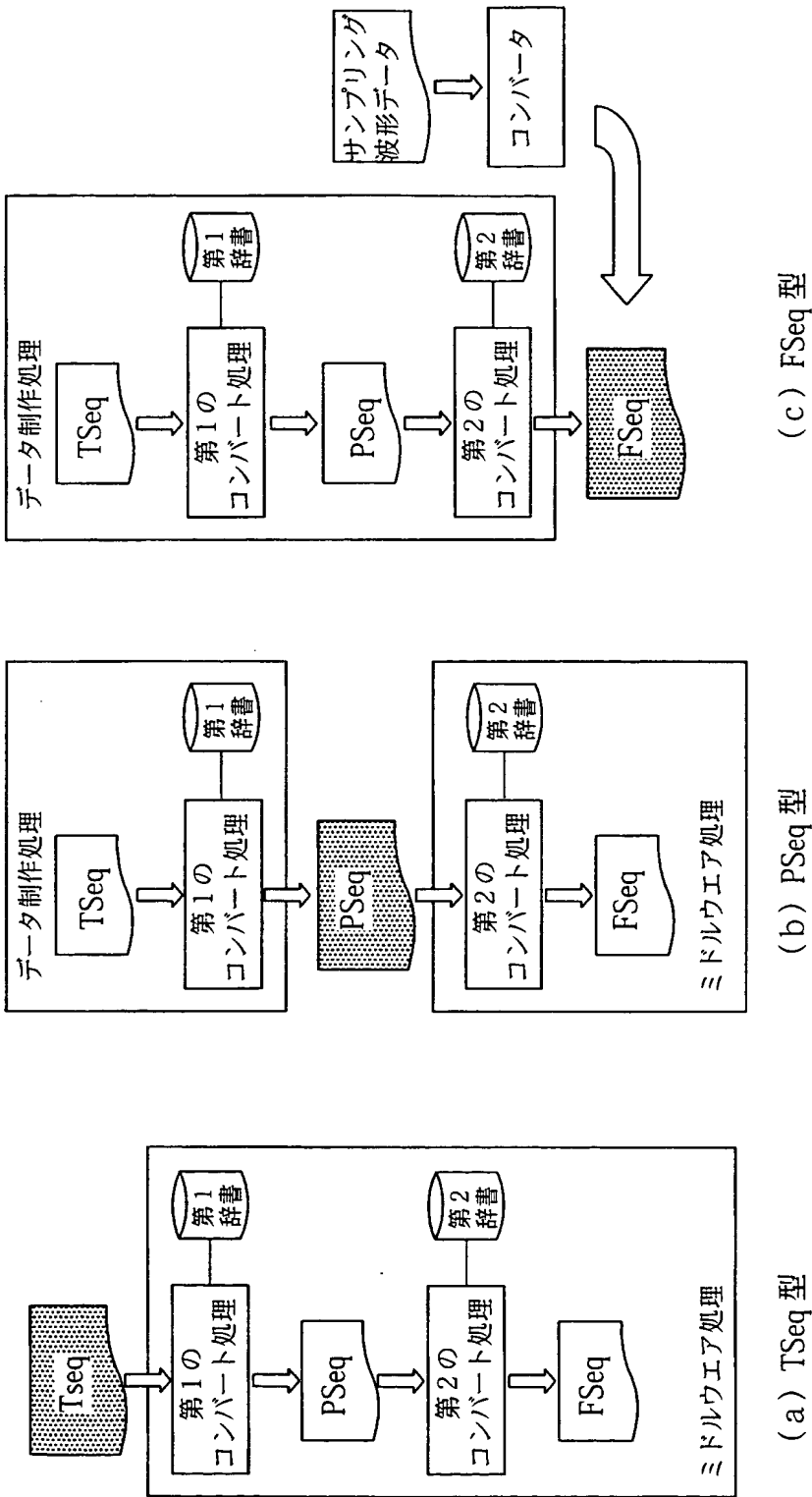
【図 3】



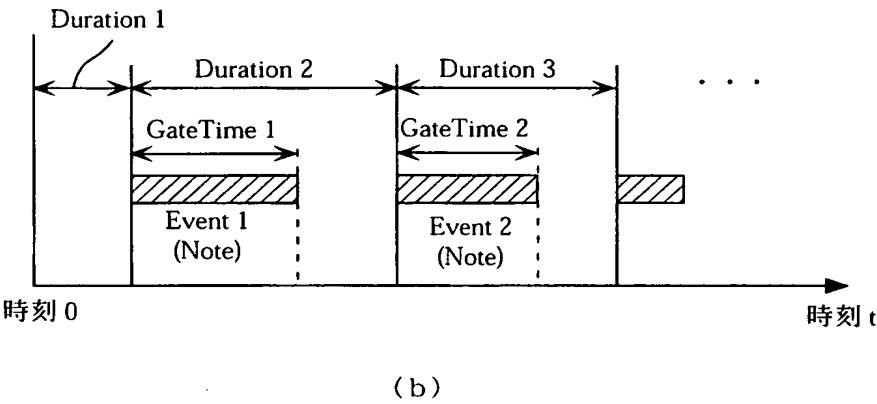
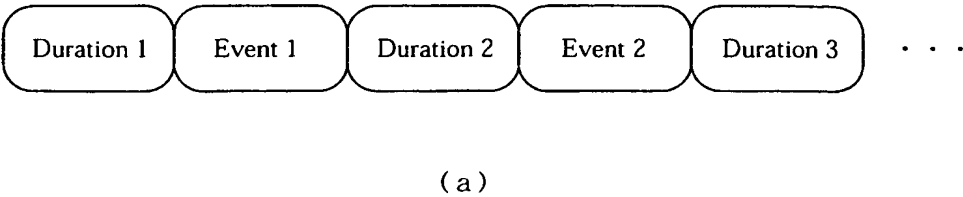
【図 4】



【図5】



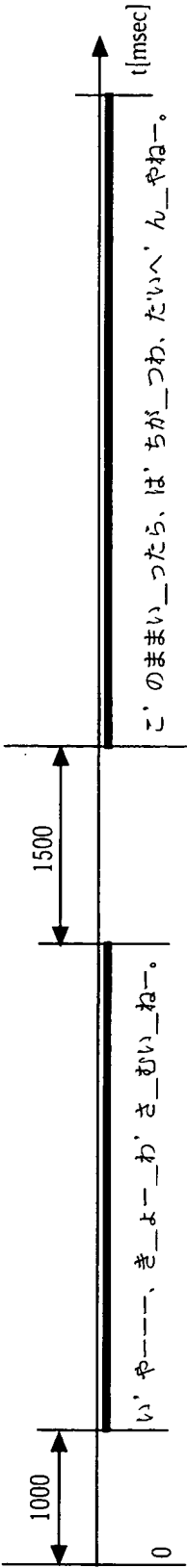
【図 6】



【図 7】

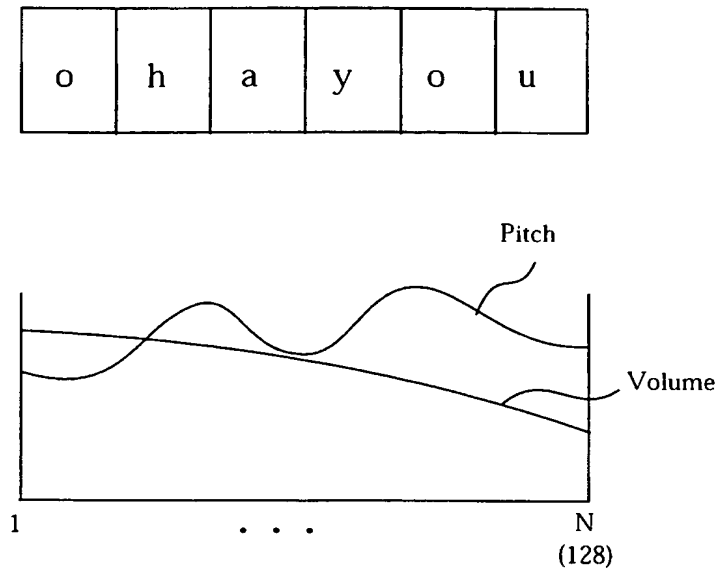
<LJAPANESE
<CS-JIS
<G4
<V1000
<N64
<PI1000
<Tい' や----、き__よ__わ' さ__むい__ね-。
<PI500
<Tこ' のままい__つたら、は' ちが__つわ、たいへ' ん__やね-。
<Q

(a) TSeq Sample

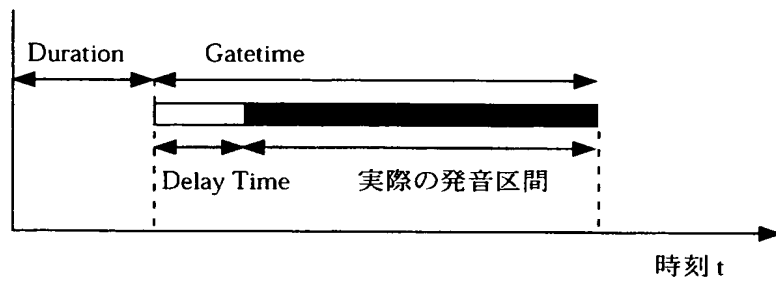


(b)

【図 8】

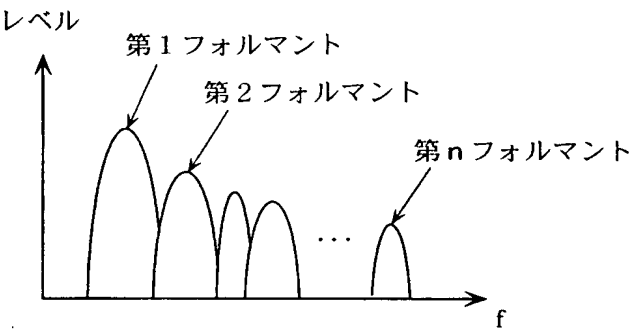


【図 9】

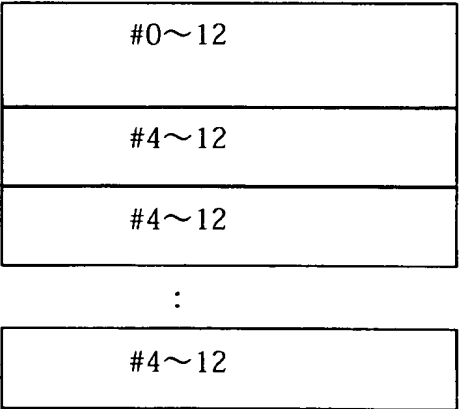


ゲートタイム長とDelay Timeとの関係

【図 1 0】

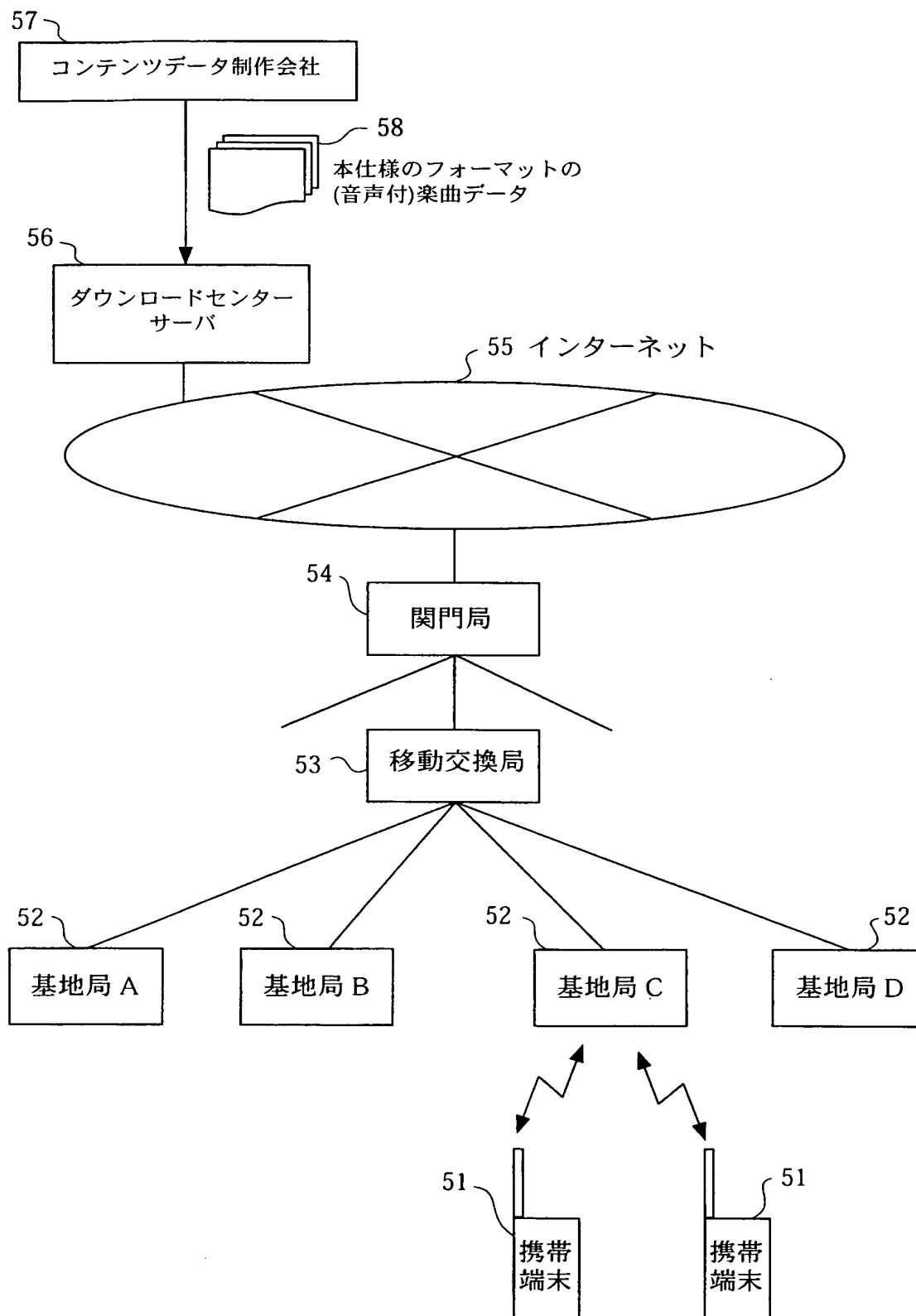


【図 1 1】

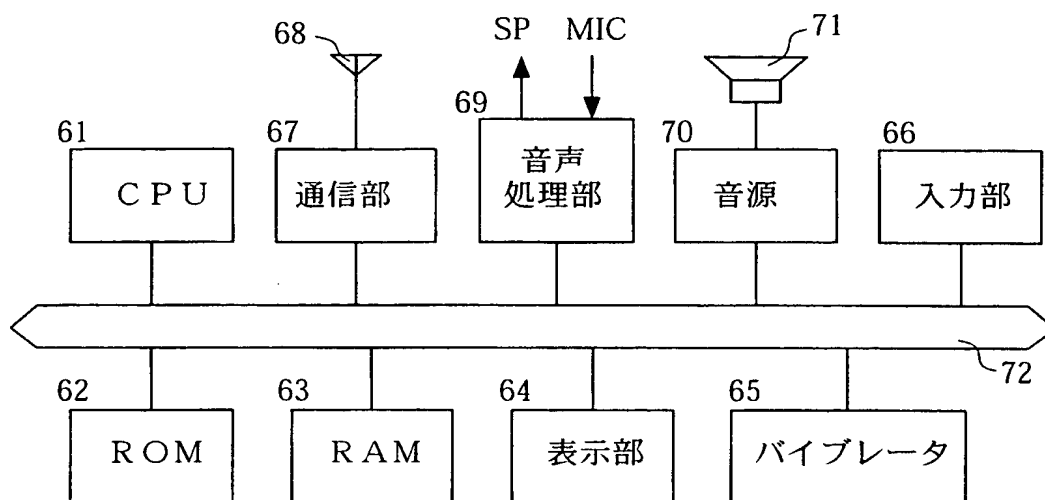


FSeq Data Chunk の Body

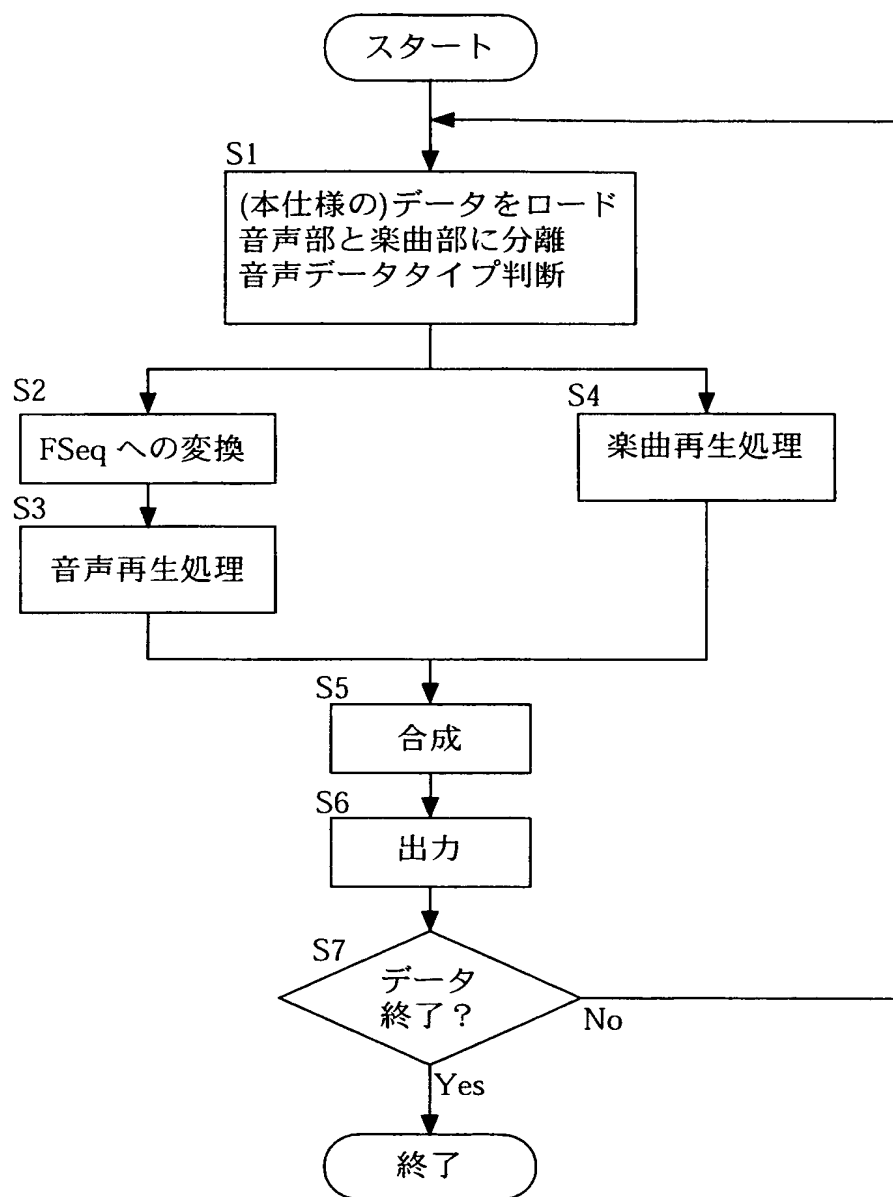
【図 12】



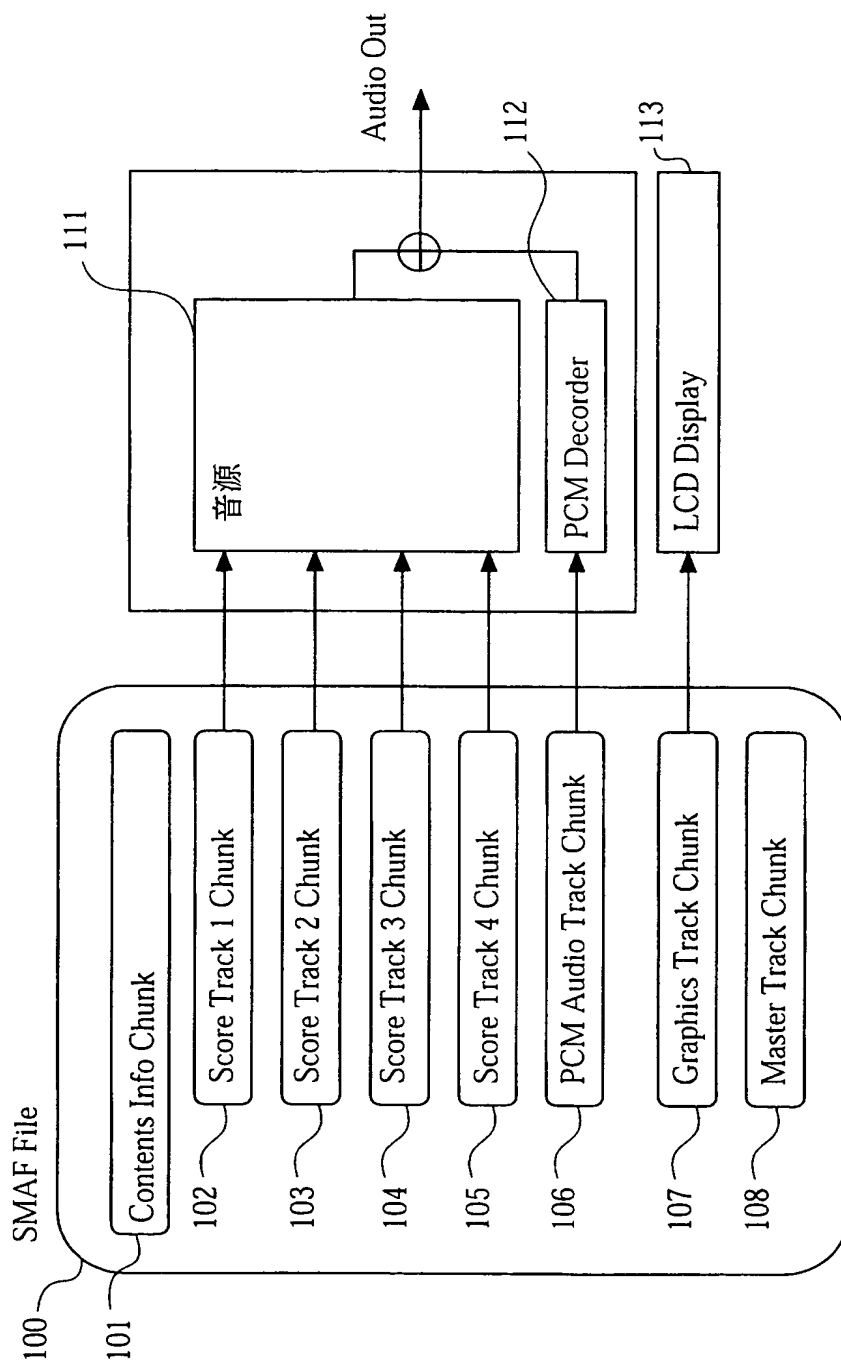
【図 13】



【図 14】



【図 15】



【書類名】 要約書

【要約】

【課題】 楽曲シーケンスデータと音声再生シーケンスデータを同期して再生することができるシーケンスデータ交換フォーマットを提供する。

【解決手段】 ファイル 1 はチャンク構造となっており、内部に管理用の情報を含むコンテンツ・インフォ・チャンク 2、オプション・データ・チャンク 3 及び音声再生用の H V トラックチャンク 4 を含む。H V トラックチャンク 4 に含まれる音声再生シーケンスデータは、（１）合成される音声の読みを示すテキスト情報と音声表現を指定する韻律記号とからなるテキスト記述型、（２）合成される音声を示す音素情報と韻律制御情報とからなる音素記述型、又は、（３）再生される音声を示すフレーム時間毎のフォルマント制御情報からなるフォルマントフレーム記述型のいずれかを選択することができる。H V トラックチャンク 4 は、SMA F ファイル中にスコアトラックチャンクなどと同様に含ませることができる。

【選択図】 図 1

特願 2 0 0 2 - 3 3 5 2 3 3

出 願 人 履 歴 情 報

識別番号

[0 0 0 0 0 4 0 7 5]

1. 変更年月日

1 9 9 0 年 8 月 2 2 日

[変更理由]

新規登録

住 所

静岡県浜松市中沢町 1 0 番 1 号

氏 名

ヤマハ株式会社